# Architecture & Design Characteristics/Controls

NIST Big Data Working Group, Definitions and Taxonomy Subgroup

UCSD, Super Computing Facility
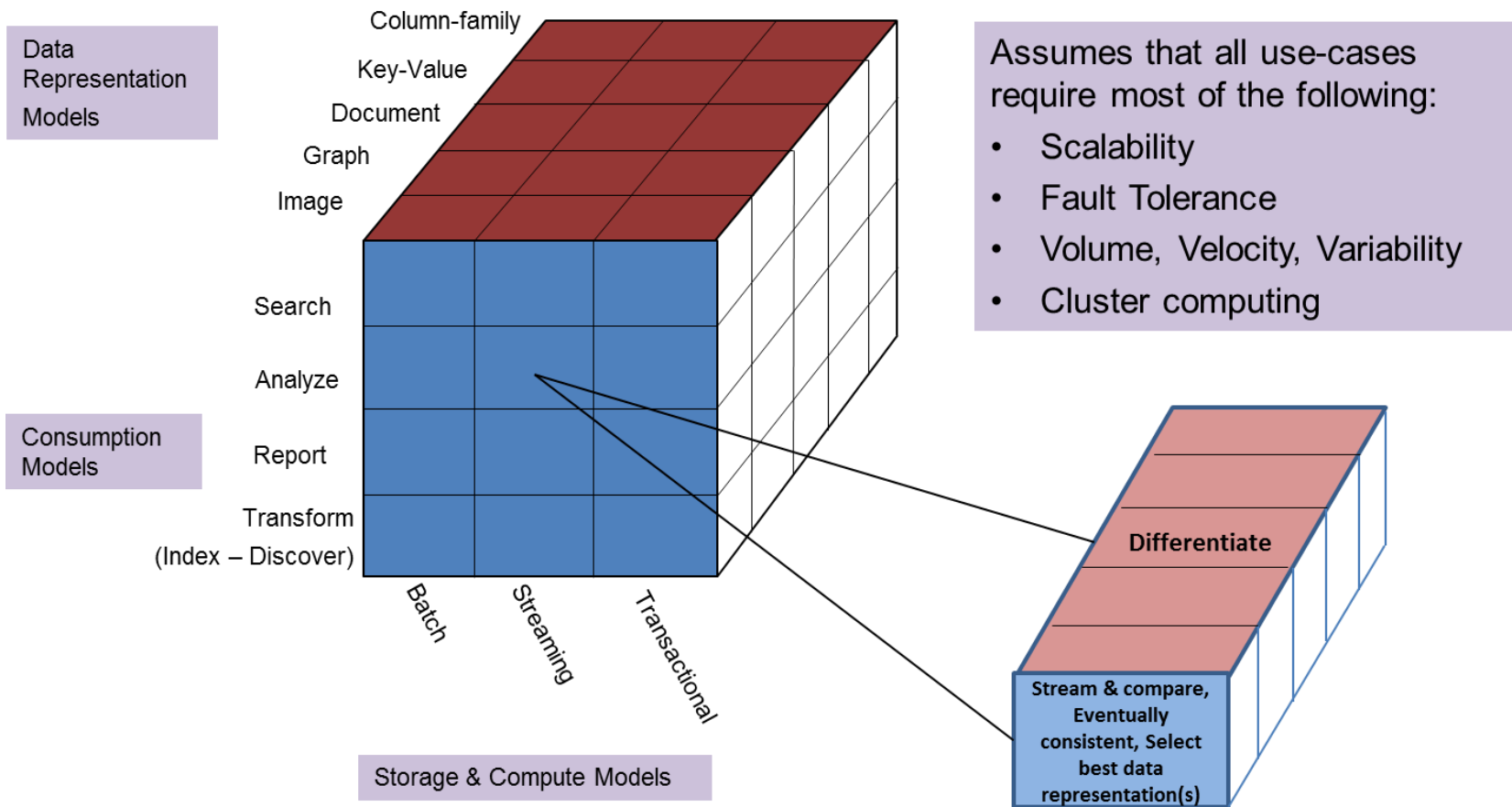March 18-21, 2012

John Rogers, Master Technical Consultant, HP
Nancy Grady, PhD, Technical Fellow, SAIC

# Big Data Architecture Characteristics Cube



Data Representation Models

Consumption Models

Storage & Compute Models

Column-family
Key-Value
Document
Graph
Image
Search
Analyze
Report
Transform (Index – Discover)

Batch
Streaming
Transactional

Assumes that all use-cases require most of the following:

- Scalability
- Fault Tolerance
- Volume, Velocity, Variability
- Cluster computing

Differentiate

Stream & compare, Eventually consistent, Select best data representation(s)

# Big Data Projects Risk Avoidance?



Business Purpose, Concerns, Risks, Guidance, and Execution

# by way of Example

**NIST Special Publication 800-53**
Revision 4

---

## Security and Privacy Controls for Federal Information Systems and Organizations

# Our big data assignment, today

Provide detailed guidance to the solutions architects and designers to:

- Build me a system that streams web page hits to a classification model and spits out alerts if a customer meets or exceeds a high-interest thresh-hold.
- Provide continuous monitoring and validation of algorithm performance.
- Provenance is unimportant.
- Data consumers do not need special analysis, fusion, or visualization tools.
- This is primarily an alerting system.
- Scale is 2TB's per week.  Retain history for 5 years. Plan for future expansion.
- Provide assurances that the system will work as planned

# How do we approach this?

By analogy

**Follow the pattern of NIST SP 800-53, Version 4, 4/30/2013**

**"Security and Privacy Controls for Federal Information Systems and Organizations"**

**http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-53r4.pdf**

TABLE 1: SECURITY CONTROL IDENTIFIERS AND FAMILY NAMES

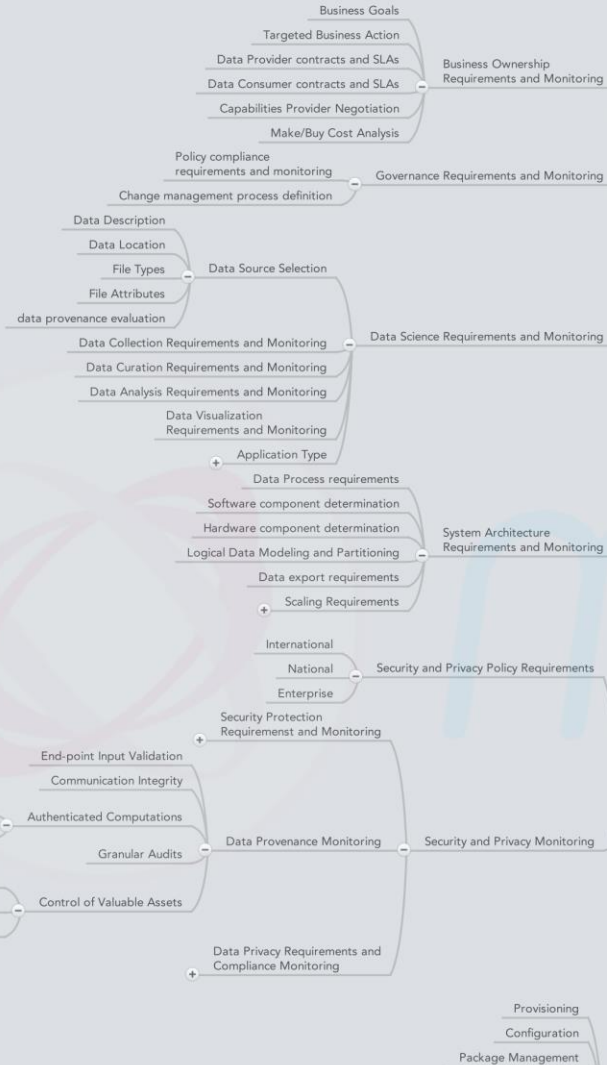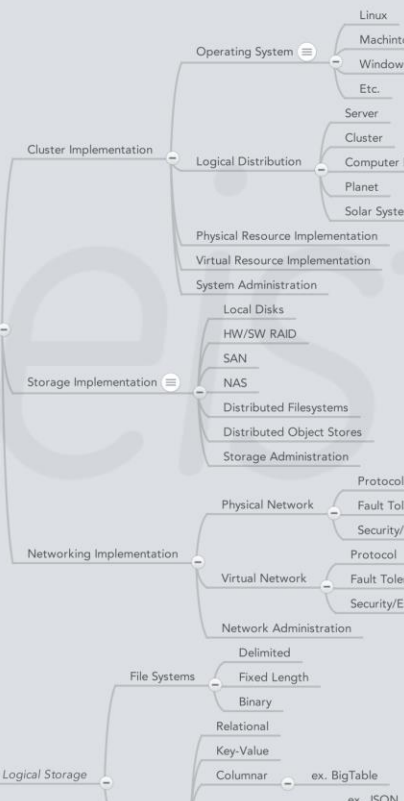| ID | FAMILY | ID | FAMILY |
|----|--------|----|--------|
| AC | Access Control | MP | Media Protection |
| AT | Awareness and Training | PE | Physical and Environmental Protection |
| AU | Audit and Accountability | PL | Planning |
| CA | Security Assessment and Authorization | PS | Personnel Security |
| CM | Configuration Management | RA | Risk Assessment |
| CP | Contingency Planning | SA | System and Services Acquisition |
| IA | Identification and Authentication | SC | System and Communications Protection |
| IR | Incident Response | SI | System and Information Integrity |
| MA | Maintenance | PM | Program Management |

Big Data Taxonomies
- Level 1: Roles
- Level 2: Activities
- Level 3: Components
- Level 4: Sub Components

# Big Data Taxonomies

| Architecture Characteristics (Controls) | | | | | |
|---|---|---|---|---|---|
| **Business Ownership & Monitoring** | **Data Science Requirements & Monitoring** | **Security & Privacy Policy** | **System Management** | **Framework Provider** | **Data Application Provider** |
| Goals | Collection | International | In-House | Platform | Collection |
| Targeted Action | Curation | National | Data Center | Processing | Preparation |
| Data related SLAs | Analysis | Enterprise | Cloud | Infrastructures | Analytics |
| Capabilities Provider Negotiation | Data Visualization | Monitoring | | | Visualization |
| Make/Buy Analysis | Application Type | Auditing | | | Access |

# Characteristic (Control) Identifiers and Family names

| ID | Family | ID | Family |
|----|--------|----|--------|
| BO | Business Ownership Monitoring | DC | Data Consumer |
| DS | Data Science Requirements Monitoring | DP | Data Provider |
| SP | Security & Privacy  Policy | **FP** | **Framework Provider** |
| SM | System Management | DA | Data Applications Provider |

Family Names are consistent with Level 1 (Roles) in the Big Data Taxonomy

# Example: Drill-down into Framework Provider activities

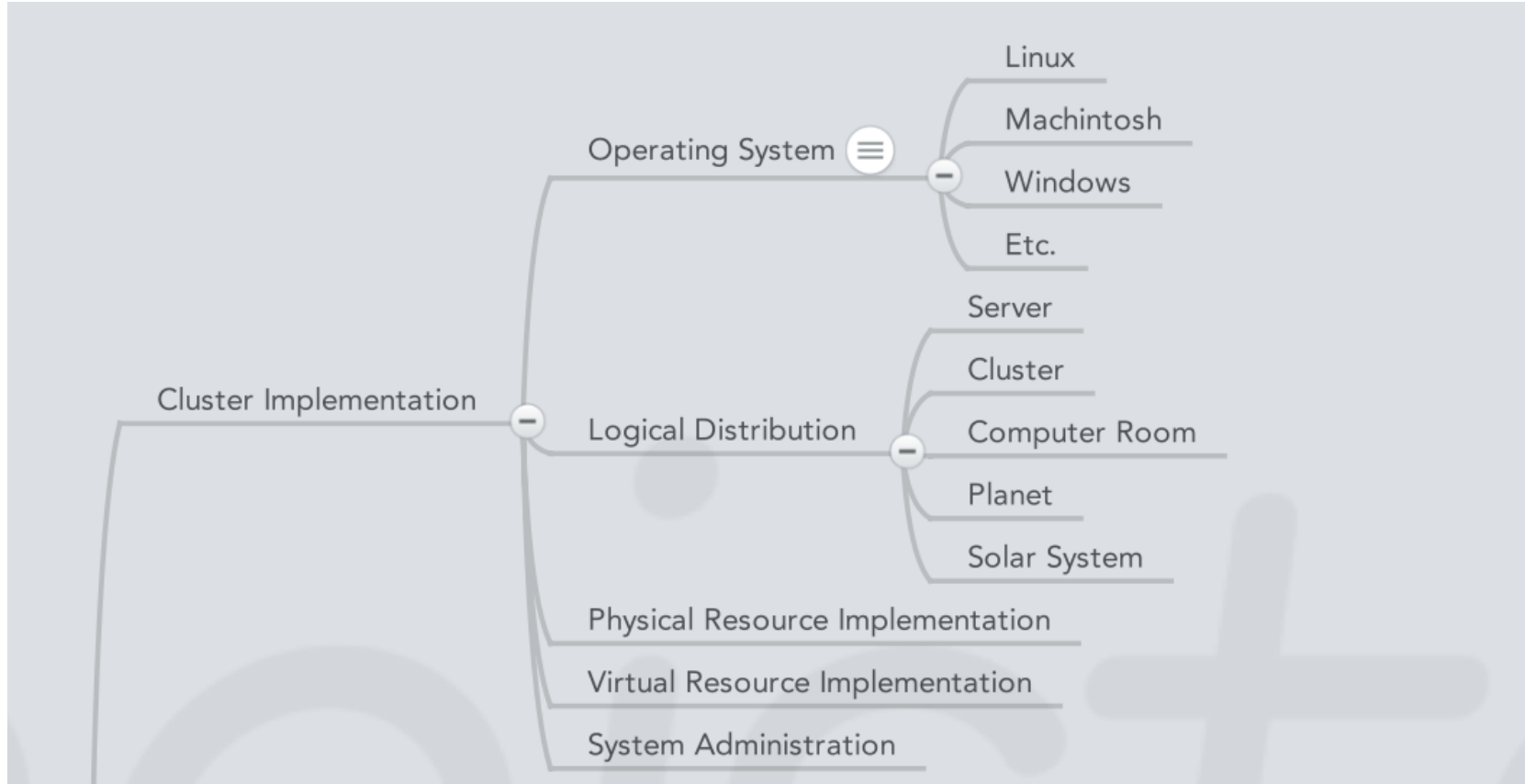| Framework Provider Family (FP) | | | |
|---|---|---|---|
| ID | Activity | | |
| **IN** | **Infrastructure** | | |
| PL | Platform | | |
| PF | Processing Framework | | |

Activity categories are consistent with Level 2 (Activities) in the Big Data Taxonomy

# Drill down to Framework Provider, **Infrastructure Activity components**

| Framework Provider Family (FP) | | |
|---|---|---|
| Infrastructure Activity (I) | | |
| ID | Component | |
| **CL** | **Cluster** | |
| **ST** | **Storage** | |
| NW | Network | |

Implementation categories are consistent with Level 3 (Components) in the Big Data Taxonomy

# Infrastructure Cluster Implementation

Cluster Implementation
- Operating System
  - Linux
  - Machintosh
  - Windows
  - Etc.
- Logical Distribution
  - Server
  - Cluster
  - Computer Room
  - Planet
  - Solar System
- Physical Resource Implementation
- Virtual Resource Implementation
- System Administration

# Drill down to Framework Provider, Infrastructure Activity, Cluster and Storage Sub-components

| Framework Provider Family (FP) | | |
|---|---|---|
| Infrastructure Activity (I) | | |

| Cluster (C) | | Storage (S) | |
|---|---|---|---|
| ID | Sub-Component | ID | Sub-Component |
| DR | Logical Distribution | **DF** | **Distributed File System** |
| **PR** | **Physical Resource** | DO | Distributed Object Store |
| VR | Virtual Resource | RD | RAID |
| OS | Operating System | **DD** | Disk type **(HDD, SDD,** Array, Network) |
| SA | System Administration | SN | Storage Administration |

Sub-categories are consistent with Level 4 (Sub-Components)

# Drill down to Framework Provider, Infrastructure Activity, Cluster Sub-component, Physical Resource

| Framework Provider Family (FP) | | |
|---|---|---|
| Infrastructure Activity (I) | | |
| Cluster  Component (C) | | |
| Physical Resource Sub-component (PR) | | |
| ID | Sub-Component | Commodity Server Performance | |
| 1 | **Commodity Server** | 1 | Base |
| 2 | Server-SAN | **2** | **Mid** |
| 3 | Custom Server | 3 | High |

**FP-I-C-PR-12 (Mid-range server)**

# Additional Characteristics (controls)

Guidance for today's assignment

FP-I-C-PR-13 (Mid Range Server)

BO-MB (High)

SP-PM (low)

DS-CM (Med)

DS-AM (Med)

DC-VA (High)

DC-SR (Low)

DP-WC (High)

DA-AP-ML-CL (High)

DA-AP-ML-DF (Med)

DA-AP-ML-SA (Low)

- Build me a system that *streams* web page hits in a *classification* model and spits out *alerts* if a customer meets or exceeds a high-interest thresh-hold.
- Provide continuous monitoring and validation of algorithm performance.
- Provenance is not important.
- Data consumers do not need special analysis, fusion, or visualization tools.
- This is primarily an alerting system.
- Scale is TB's per week.
- Keep capital cost reasonable.
- Infrastructure selection based on testing

# Again, by way of Example

# Ex: Evaluate Workloads

## TeraSort

- Base workload: Find, Shuffle, Sort in Order
- Good Overall utilization (CPU, Memory, Disk, Network IO)
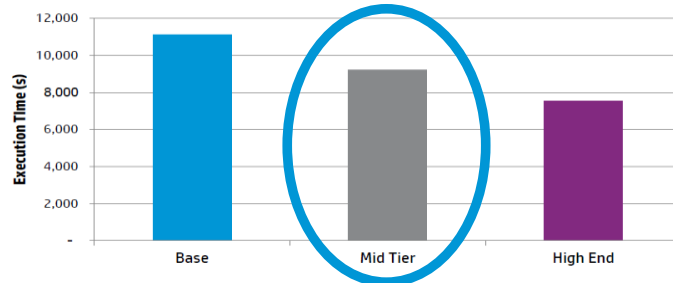
## Intel Hi-Bench Workloads

- Web-Search: K-V Indexing
- Machine Learning: K-Means Clustering

## Data Analytics

- Query on data warehouse type data
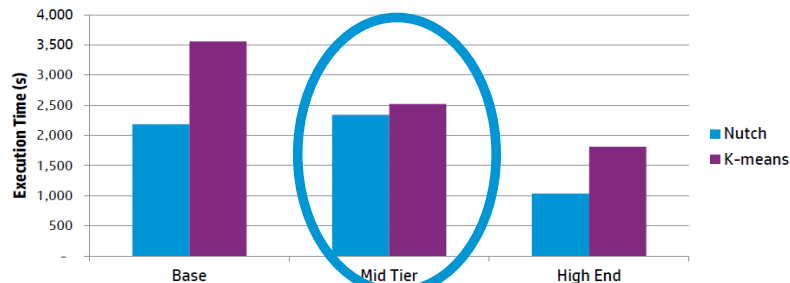- Complex queries with many joins and grouping/sorting operations

**Note: Tests in 2013 using CDH 3.x**



**TeraSort**
2TB TeraSort execution time results – lower is better

**Nutch Indexing and K-means Clustering**
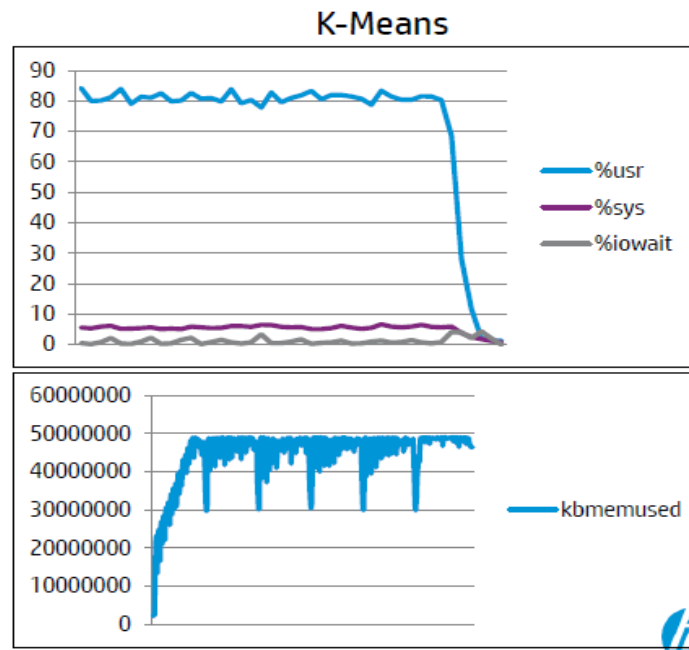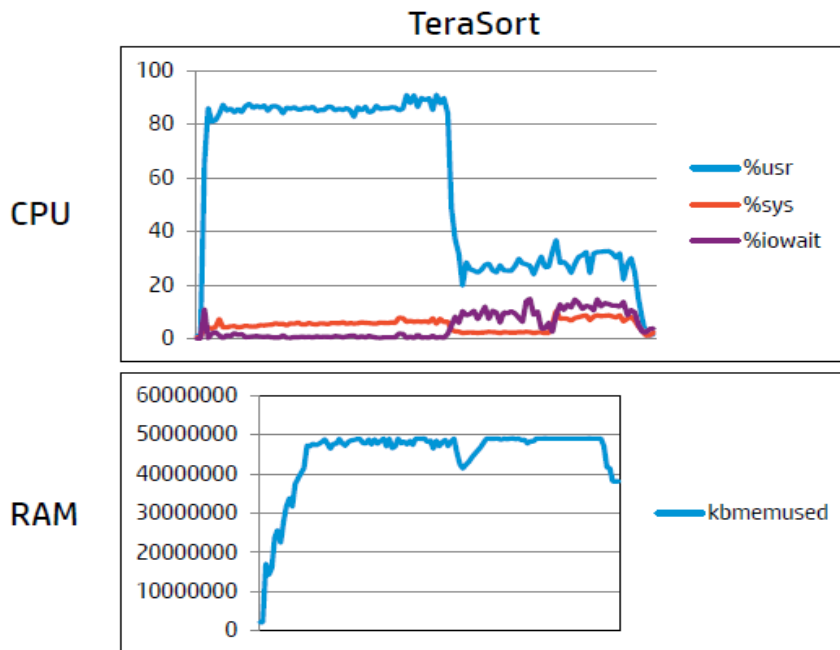Execution time in seconds – lower is better

**Data Analytics – Hive Queries**
Execution time in seconds – lower is better

# Ex: Evaluate limits



TeraSort — CPU: %usr, %sys, %iowait; RAM: kbmemused
K-Means — CPU: %usr, %sys, %iowait; RAM: kbmemused

- Typically CPU, Memory, and or Disk I/O limited.
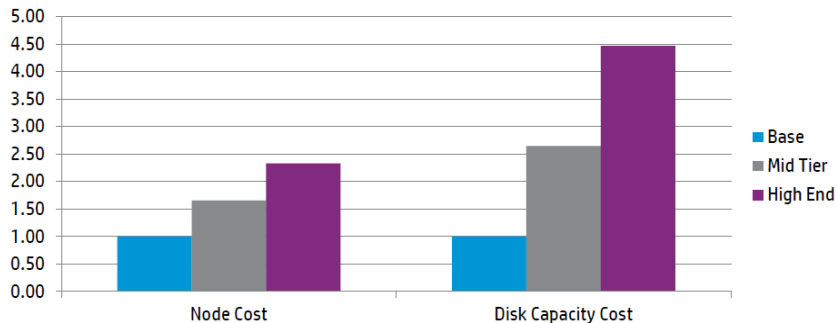- Only networked limited about 10-15% of the time, if at all

# Ex: Evaluate tradeoffs

## Performance vs. cost

- Sorting and Cluster workloads improved by 25% over base
- Node costs increase by 50% for each step
- Tradeoff improved disk reads at higher cost
- Idle power consumption improves energy use
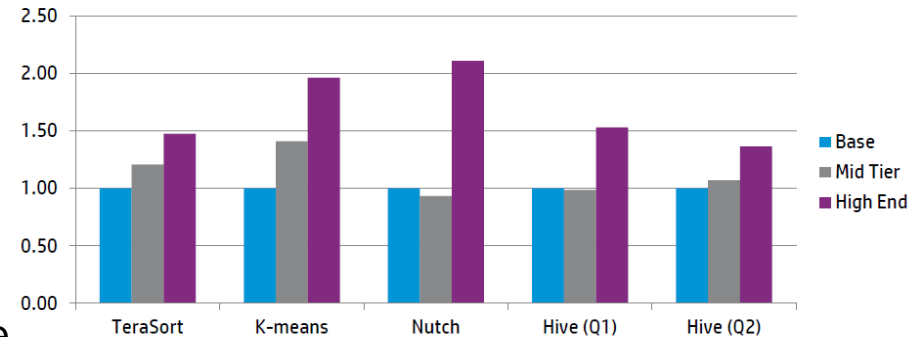- Tradeoff faster batch process at higher intermittent energy use

**Relative Performance**

Normalized results – Higher is better



**Relative Cost Comparisons**

Relative Data Node Cost and Disk Capacity Cost – Lower is better



**System Power Consumption**

Worker Node Power Consumption (Watt) – Lower is better

# Next Steps for Research

FP-I-C-PR-13 (Mid Range Server)

BO-MB (High)

SP-PM (low)

DS-CM (Med)

DS-AM (Med)

DC-VA (High)

DC-SR (Low)

DP-WC (High)

DA-AP-ML-CL (High)

DA-AP-ML-DF (Med)

DA-AP-ML-SA (Low)

# Additional Cataloguing Considerations

Between and Inside Components in the Taxonomy

- **Definitions and Taxonomy group named functional components**
- **We have not yet addressed the different ways the components work**
  - Or the different ways end-to-end systems work
- **Use cases give data-process lifecycles**
  - See Bob Marcus' M0297 high-level scenarios for use case categorization
- **Need to figure out the dimensions that differentiate component instantiations**
  - e.g. Inter-node communication -> implying latency in consistency
  - e.g. data location for processing (in-memory, on disk,…)
  - e.g. fault tolerant scheme (replication, master-slave, …)
  - e.g. analytics time constraints (streaming, interactive, batch,…)

# Goal of Research

Convert Big Data WG Architecture and Use-Cases to Characteristic (Control) codes

**Present Concept**

- Create Taxonomies, Use cases, and architecture features
- Group functional components

**Research**

- Address component characteristics and behaviors
- Characterize different system/implementation behaviors
- Identify critical differentiators
- Gather published test results

**Build a List of Characteristics (Controls)**

- Map component characteristics, behaviors, and differentiators
- Assign codes
- Create evaluation criteria (low, medium, high)

**Brief the Big Data Working Group**

- Working Group briefings

**Final Research Paper**

- Publish or perish

**Your feedback is important to us. Please take a few minutes to complete the session survey.**

# Thank you

# Questions?

John.W.Rogers@hp.com
Nancy.W.Grady@saic.com